

- Residual error ε_i 's are independent but the residuals e_i 's are not since we have n observations but the degrees of freedom is only $n - k$!
- e_i 's are the observed values of ε_i
- It is convenient to think of the residuals as the values of the errors. Hence, we can test the underlying assumptions on the residuals.

One idea is to plot the residuals in an effective way to test how well the model fits the data & how well the assumptions are satisfied.

* Types of Residuals :

1) Regular Residuals (e_i) :

$$e_i = y_i - \hat{y}_i$$

$$\begin{aligned} e &= y - \hat{y} \\ &= y - X\hat{\beta} \\ &= y - X(X^T X)^{-1} X^T y \\ &= (I - X(X^T X)^{-1} X^T) y \\ &= (I - H) y \end{aligned}$$

H is symmetric and idempotent as :

$$H^2 = (X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T)^T$$

$$\begin{aligned}
&= (X(X^T X)^{-1} X^T) (X(X^T X)^{-1} X^T) \\
&= X(X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T \quad \text{switch the order!} \\
&= X(X^T X)^{-1} X^T \\
&= H
\end{aligned}$$

H is called the Hat Matrix.

$$\begin{aligned}
\text{Hence, } e &= (I - H) Y \\
&= (I - H)(X\beta + \varepsilon) \\
&= X\beta - HX\beta + (I - H)\varepsilon \\
&= X\beta - X(X^T X)^{-1} X^T X\beta + (I - H)\varepsilon \\
&= X\beta - X\beta + (I - H)\varepsilon \\
&= (I - H)\varepsilon
\end{aligned}$$

$$\begin{aligned}
\text{Var}(e) &= \text{Var}((I - H)\varepsilon) \\
&= (I - H)^2 \text{Var}(\varepsilon) \\
&= (I - 2H + H^2) \sigma^2 \\
&= (I - 2H + H) \sigma^2 \\
&= (I - H) \sigma^2
\end{aligned}$$

$$= \begin{pmatrix} \text{Var}(e_1) & \text{Cov}(e_1, e_2) & \dots & \text{Cov}(e_1, e_n) \\ \text{Cov}(e_2, e_1) & \text{Var}(e_2) & \dots & \text{Cov}(e_2, e_n) \\ \vdots & \vdots & \ddots & \vdots \end{pmatrix}$$

$$\begin{aligned}
\text{Var}(e_i) &= (1 - h_{ii}) \sigma^2 & h_{ii} &= u_i(X^T X)^{-1} u_i \\
\text{Cov}(e_i, e_j) &= (\rho - h_{ij}) \sigma^2 = -h_{ij} \sigma^2 & I_{i,j} &= \rho \text{ (diagonal)}
\end{aligned}$$

h_{ii} measures the distance of the i^{th} observation from the center at the n coordinate. $0 \leq h_{ii} \leq 1$.

b) Studentized Residuals (r_i) :

$$\text{Var}(e_i) = \sigma^2 (1 - h_{ii}) = \text{MSRes} (1 - h_{ii})$$

$$\mathbb{E}(e_i) = 0$$

$$r_i = \frac{e_i}{\sqrt{\text{MSRes} (1 - h_{ii})}}$$

studentized residuals have a constant variance
 $\text{var}(r_i) = 1$ regardless of the location of the n -
coordinates when the form of the model is correct.

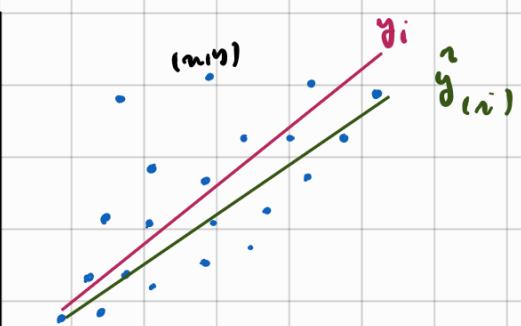
c) Standardized Residuals (d_i) :

$$\text{Var}(e_i) \approx \sigma^2$$

$$d_i = \frac{e_i}{\sqrt{\text{MSRes}}}$$

d) Press Residuals ($e_{(i)}$)

$$i^{\text{th}} \text{ press residual : } e_{(i)} = y_i - \hat{y}_{(i)}, \quad e_i \neq e_{(i)}$$



- predicted line with all the observations
- predicted line with the

This residual is useful in the presence of an outlier.
 \hat{y}_i is the fitted value at the i^{th} response based on all the observed except the i^{th} .

e.g.: $M_1 \rightarrow Y = 2x + 3$ $M = (2, 4)$
 $M_2 \rightarrow Y = 5x + 9$

$$\begin{array}{ll} \hat{Y}_{M_1} = 7 & Y - \hat{Y}_{M_1} = 4 - 7 = -3 \\ \hat{Y}_{M_2} = 19 & Y - \hat{Y}_{M_2} = 4 - 19 = -15 \# \end{array}$$

- We delete the i^{th} observation then fit the response model to the remaining $(n-1)$ observations and predict y_i .

- $e_{(i)} = y_i - \hat{y}_{(i)} = \frac{e_i}{1 - h_{ii}}$

- If we have very large press residual, we can identify the observations where the model does not fit the data well.

- Press Statistic $= \sum_{i=1}^n e_{(i)}^2$
 $= \sum_{i=1}^n (y_i - \hat{y}_{(i)})^2$
 $= \sum_{i=1}^n \left(\frac{e_i}{1 - h_{ii}} \right)^2$

which measures how well a regression model will perform in predicting a new data point.

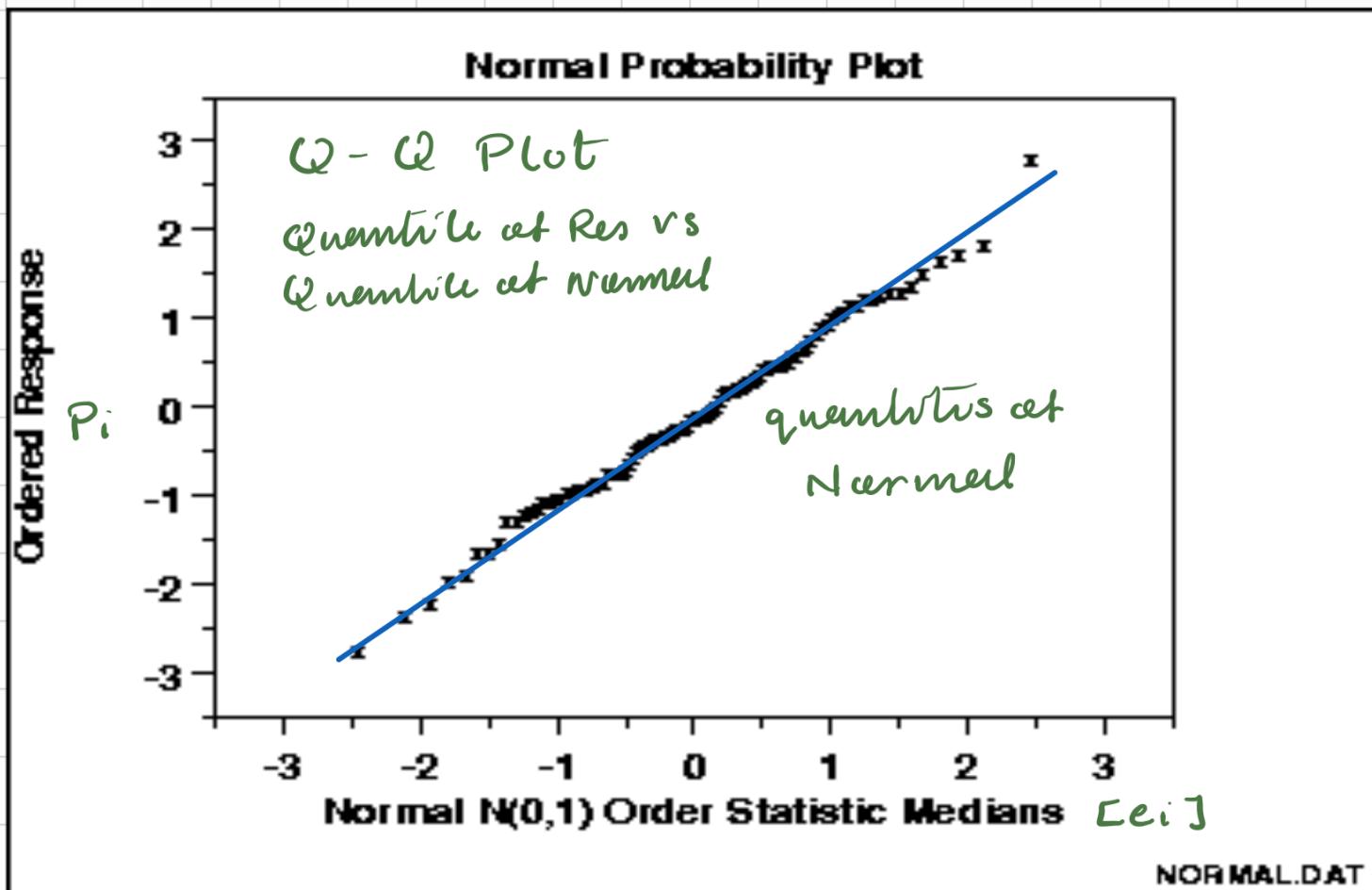
* Types of Residual Plots :

1) Normality Plot :

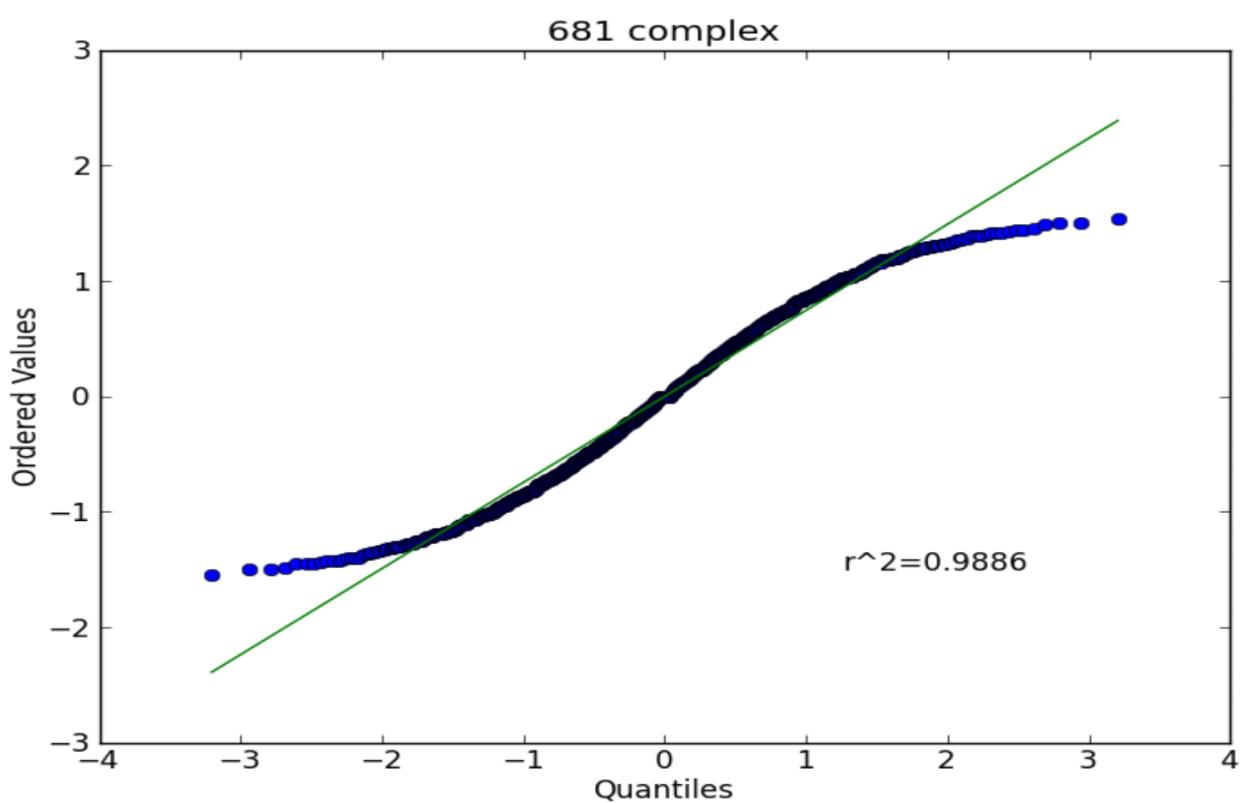
Let e_1, \dots, e_n be the n residuals, $e_{(1)}, \dots, e_{(n)}$ are the residuals in an increasing order of magnitude.

Plot $[e_i]$ against Cumulative Probabilities :

$$P_i = \frac{i - \gamma_2}{n}$$

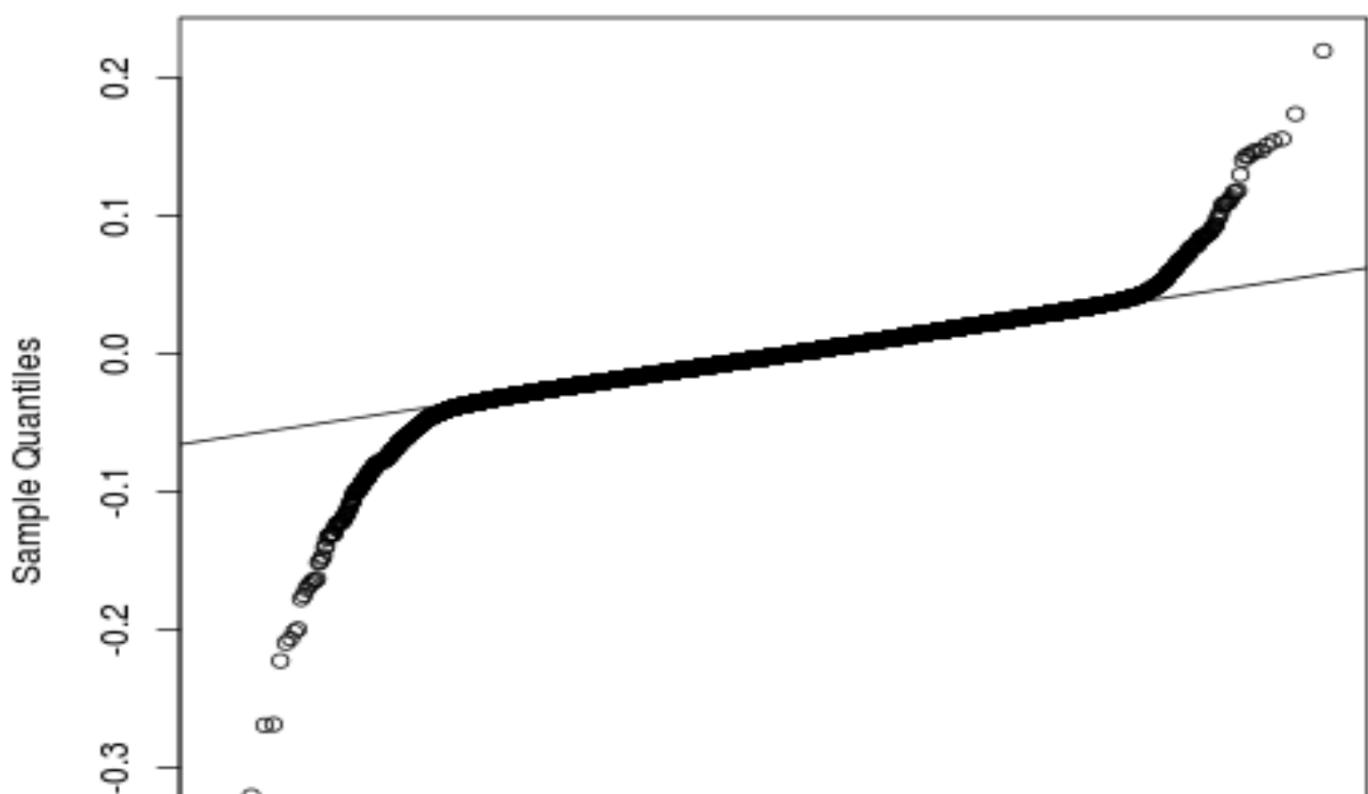


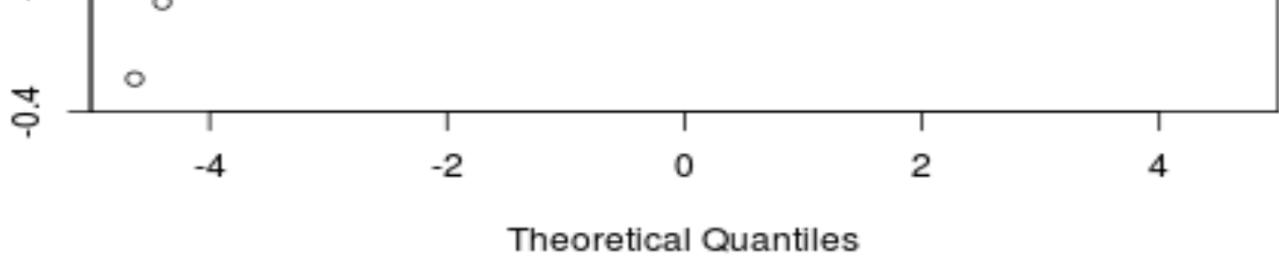
Normal Distribution (Ideal Situation)



Light Tailed Distributions

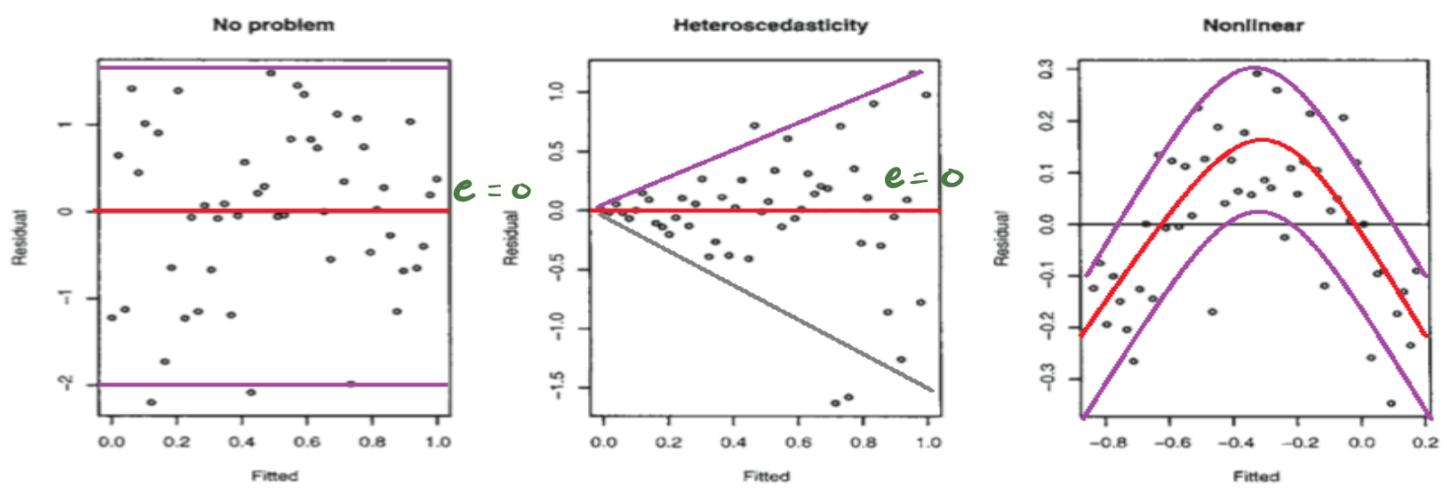
Normal Q-Q Plot





Heavy Tailed Distributions

a) Plot of Residuals e_i vs fitted values (\hat{y}_i)



Satisfactory
Model

Unsatisfactory
Model

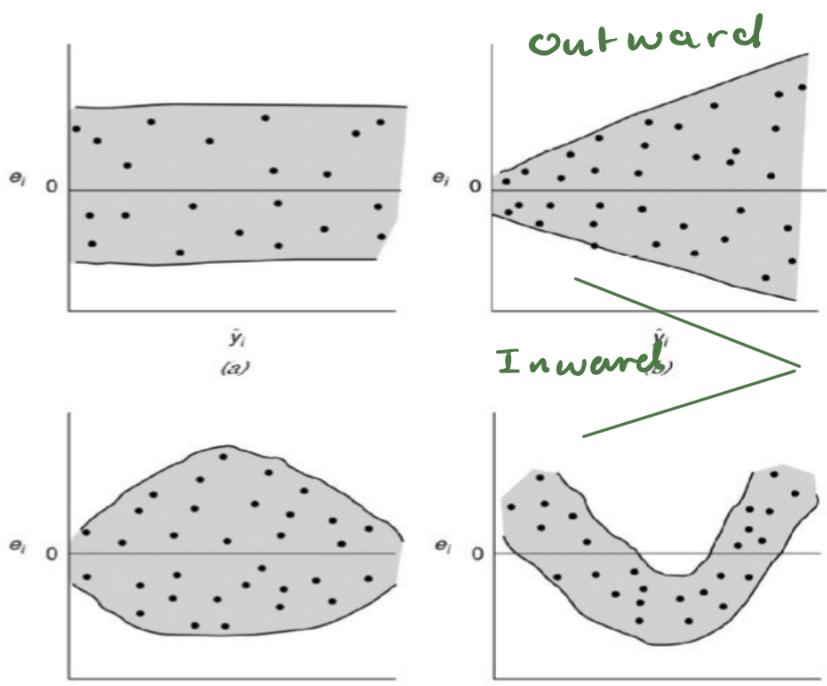


Figure 4.3 Patterns for residual plots: a) satisfactory; b) funnel; c) double bow; d) nonlinear.

a) Good regression model :

scatter in residuals is roughly constant wrt \hat{y}_i and centered around $c = 0 \therefore \text{Var}(\varepsilon) = \sigma^2$ is constant variance.

b) Outward Opening funnel : $\text{Var} \varepsilon \nearrow$ as $\hat{y}_i \nearrow$
Inward Opening funnel : $\text{Var} \varepsilon \searrow$ as $\hat{y}_i \nearrow$

The variance $V(\varepsilon) \neq \sigma^2$ non constant

c) Dumbbell Box Structure :

Indicates a non-constant variance $V(\varepsilon) \neq \sigma^2$ happens when y is proportioned $0 \leq y \leq 1$.

d) Non-Linear Structure :

Other regression variables are needed in the model, consider some extra terms (square term x^2) to the model or apply a transformation to y (log or sqrt).

Why do we plot Res against \hat{y}_i or n but not y ?

We assume e_i and y_i are correlated but e_i and \hat{y}_i are not correlated.

Assume there exists a regression model;

$$e_i = \beta_0 + \beta_1 y_i + \varepsilon_i \quad \text{SLR}$$

$$\text{LSE} \Rightarrow \hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} = \frac{s_{ey}}{s_{yy}} = \frac{\sum (e_i - \bar{e})(y_i - \bar{y})}{\sum (y_i - \bar{y})^2}$$

$$s(e, Y) = \frac{s_{ey}}{s_{yy}} = \frac{\sum e_i (y - \bar{y})}{SST}, \text{ as } \bar{e} = 0$$

$$\begin{aligned} \bar{e} \sum (y_i - \bar{y}) &= 0 = \frac{\sum (e_i - \bar{e})(y_i - \bar{y})}{\sum (y_i - \bar{y})^2} = \frac{\sum e_i y_i}{SST} \\ \bar{y} \sum e_i &= 0 \quad \left(\frac{\sum e_i y_i}{\sum (y_i - \bar{y})^2} \right) = \frac{\sum e_i y_i}{SST} \\ &= \frac{\sum e_i y_i}{SST} = \frac{Y^T e}{SST} \end{aligned}$$

$$= \frac{Y^T (I - H) Y}{SST}$$

$$(I - H)^2 = (I - H) \quad = \frac{Y^T (I - H)(I - H) Y}{SST}$$

idempotent

$$= \frac{Y^T (I - H)(I - H) Y}{SST} = \frac{e^T e}{SST}$$

$$= \frac{e^T e}{SST} = \frac{SS_{\text{Res}}}{SST}$$

$$= \frac{\sum e_i^2}{SST} = 1 - \frac{SS_{\text{Reg}}}{SST}$$

$$= \frac{SS_{\text{Res}}}{SST} = 1 - R^2$$

$$= 1 - R^2$$

Assume there exists a regression model;

$$e_i = \beta_0 + \beta_1 \hat{y}_i + \varepsilon_i \quad \text{SLR}$$

$$LSE \Rightarrow \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{S_{e\hat{y}}}{S_{\hat{y}\hat{y}}} = \frac{\sum e_i(\bar{e})(\hat{y}_i - \bar{\hat{y}})}{\sum (\hat{y}_i - \bar{\hat{y}})^2}$$

$$\begin{aligned} \sum e_i(\hat{y}_i - \bar{\hat{y}}) &= \sum e_i(\hat{y}_i - \bar{y}) \sum e_i \\ &= 0 - 0 \\ &= 0 \end{aligned} \quad \begin{aligned} &= \frac{\sum e_i(\hat{y}_i - \bar{\hat{y}})}{\sum (\hat{y}_i - \bar{\hat{y}})^2} \\ &= \frac{\sum e_i \hat{y}_i}{\sum (\hat{y}_i - \bar{\hat{y}})^2} \end{aligned}$$

$$\begin{aligned} \hat{y}^T e &= (x \hat{\beta})^T e \\ &= \hat{\beta}^T x^T e \\ &= ((x^T x)^{-1} x^T y)^T x^T e \\ &= y^T x (x^T x)^{-1} x^T e \\ &= y^T H (I - H) e \\ &= y^T (H - H^2) e \\ &= 0 \end{aligned} \quad \begin{aligned} &= \frac{e^T \hat{y}}{\sum (\hat{y}_i - \bar{\hat{y}})^2} \\ &= \frac{y^T (I - H) H y}{\sum (\hat{y}_i - \bar{\hat{y}})^2} \\ &= \frac{y^T (H - H^2) y}{\sum (\hat{y}_i - \bar{\hat{y}})^2} \\ &= 0 \end{aligned} \quad , \quad H = H^2$$

$$\begin{aligned} \sum e_i \hat{y}_i &= e^T \hat{y} \\ e^T \hat{y} &= ((I - H) y)^T = y^T (I - H) \\ \hat{y} &= x \hat{\beta} = x (x^T x)^{-1} x^T y = Hy \end{aligned}$$

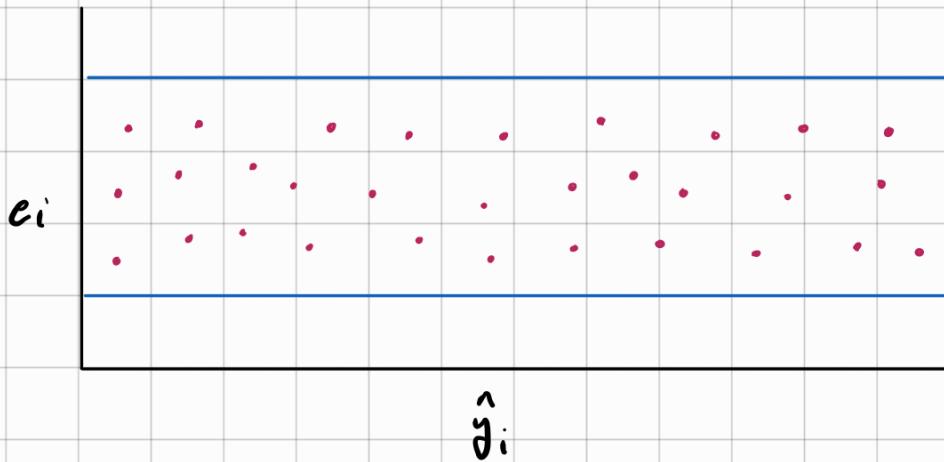
$$\begin{aligned} f(e, \hat{y}) &= 0 \Rightarrow \text{Cov}(e, \hat{y}) = \text{Cov}((I - H)y, x \hat{\beta}) = \text{Cov}((I - H)y, x (x^T x)^{-1} x^T y) \\ &= \text{Cov}((I - H)y, Hy) = (I - H)H \text{Cov}(y, y) \\ &= (H^2 - H) \text{Var}(y) = (H - H) \sigma^2 = 0 \end{aligned}$$

$\hat{\beta}_1$ for e_i & y_i : $\hat{\beta}_1 = 1 - R^2$

If $R^2 = 1$ then the slope $\hat{\beta}_1 = 0$, in e_i and y_i plot there is nothing wrong with the model as there is a theoretical justification for the relationship between the residuals and y_i , it is very likely that the residuals will not be centered within a horizontal band centered at $e = 0$, there will always be a slope

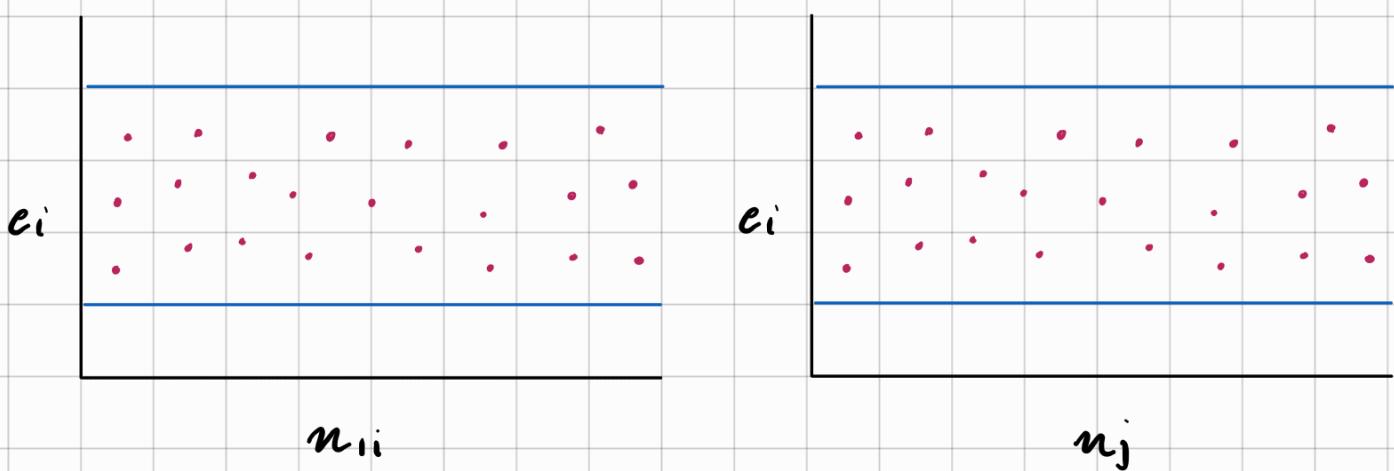
at $(1 - R^2)$ when $R^2 \neq 1$.

But when we plot e_i vs \hat{y}_i , we can get



there is no linear relationship between e_i and \hat{y}_i .

3) Partial Regression & partial residual plot



- PRP consider the marginal role of the regressor u_j given other regressors that are already present in the model (MLR case).
- In this plot, the response y_i and regressor u_j are both regressed against the other regressors in the model (except u_j) residuals are obtained from each regression.

$$f = f_1(n_1, n_2, \dots, n_{j-1}, n_j+1, \dots, n_K) \quad \& \\ n_j = f_2(n_1, n_2, \dots, n_{j-1}, n_j+1, \dots, n_K)$$

Plot of the residuals against each other gives the marginal role of the regressor n_j on y in the presence of $n_1, n_2, \dots, n_{j-1}, n_{j+1}, \dots, n_K$

Example : Consider a MLR with 2 regressors

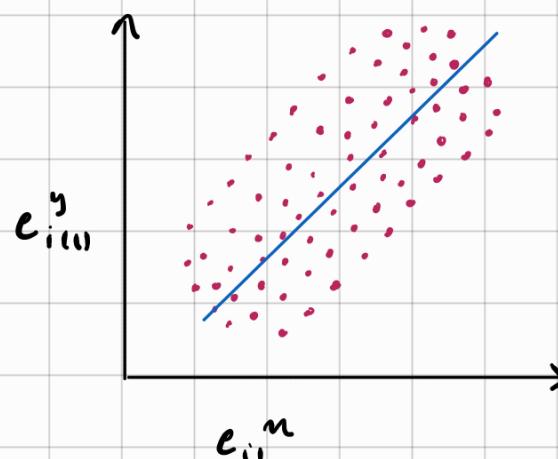
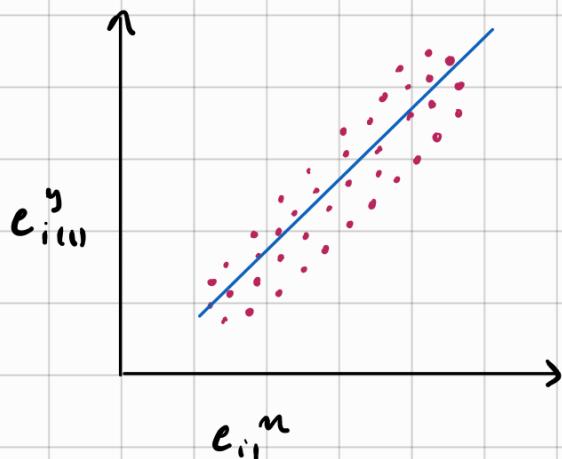
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

We are interested in the marginal role of X_1 on the response variable Y in presence of the other regressor X_2 .

$$\text{Regression of } Y \text{ on } X_2 : \hat{Y}_{i(1)} = \hat{\theta}_0 + \hat{\theta}_1 n_{i2}$$

$$\text{Regression of } X_1 \text{ on } X_2 : \hat{n}_{i(1)} = \hat{\alpha}_0 + \hat{\alpha}_1 n_{i2}$$

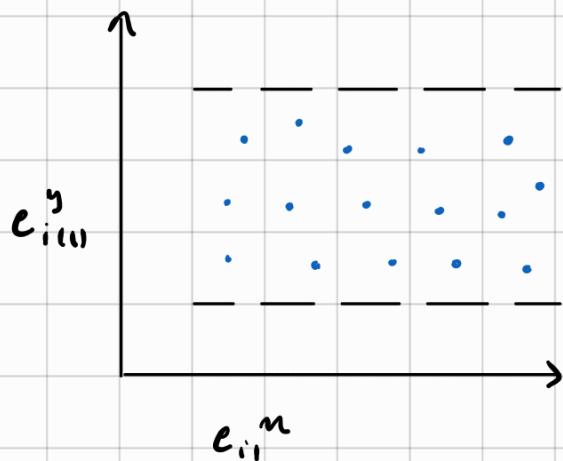
$$\text{Residuals : } e_{i(1)}^y = Y_i - \hat{Y}_{i(1)} \quad \text{eliminates the effect of } x_2 \text{ from } Y \\ e_{i(1)}^n = n_{i1} - \hat{n}_{i(1)} \quad \text{eliminates the effect of } x_2 \text{ from } x_1$$



Residual at Y has a linear relationship with the residuals at X_1 .

Adding X_1 in the model will improve the regression

framework see the marginal effect at x_i is useful as there are points in the data not explained.



In this case we do not suggest including x_i in the data but before eliminating at p_i & do some testing.

General Case : The partial residual at y for n_j is defined as $e_{i(j)}^y = y_i - \hat{y}_{i(j)}$

$\hat{y}_{i(j)}$ is the prediction based on $K-1$ regressors (not n_j)

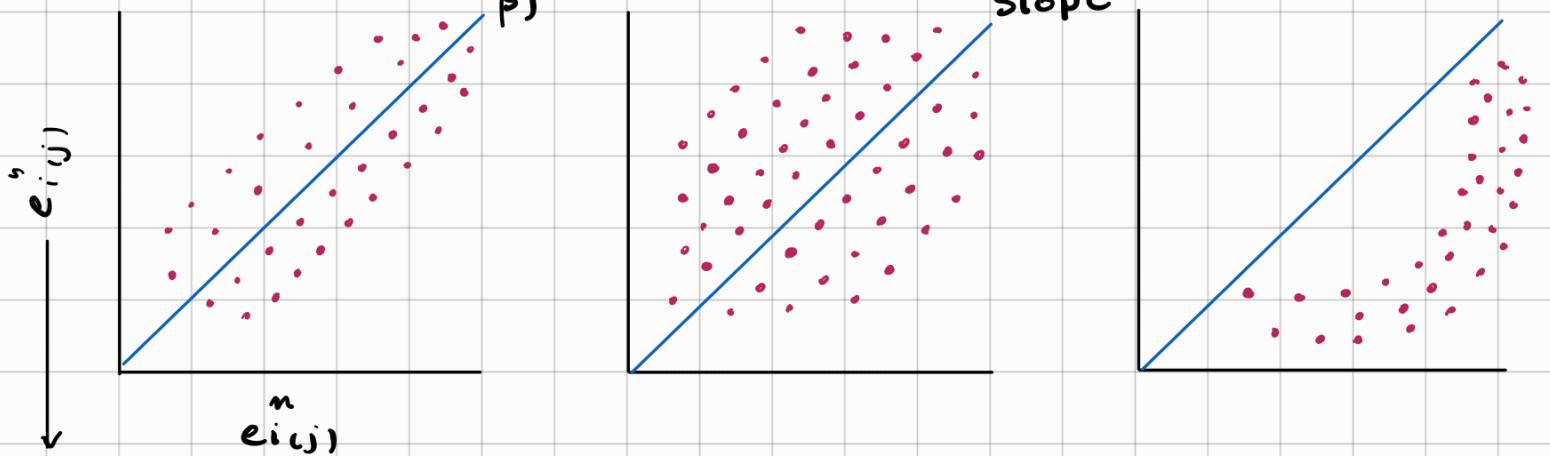
$e_{i(j)}^y$ represents the variability in y_i that isn't explained by a model that excludes the regressor n_j .

The partial residuals at n_j is defined $e_{i(j)}^n = n_{ij} - \hat{n}_{i(j)}$ where $\hat{n}_{i(j)}$ is the prediction at the regressor value n_{ij} from regression at n_j on all other regressor variables.

$e_{i(j)}^n$ is the variation that cannot be explained by $(K-1)$ other regressors .

$$e_{i(j)}^y = \beta_j e_{i(j)}^n + \varepsilon_i^*$$

β_j is the slope at the partial residual plots .



partial residual at y for n_j

less scattered plot
indicating a strong
relationship between
 y and x_j

Points are scattered

Curvilinear Bound

x_j is not linearly
related to y either
get a higher term at
 x_j or transform it
($\sqrt{x_j}$ or $\log(x_j)$)

* Checking Normality Assumption :

If not linear LSE and all the testings we know do not hold and we have to change our method of estimation

Q-Q Plot : Graphical tool used to assess normality.

It plots the sample quantiles against the theoretical quantiles on the horizontal quantile where the original data point n_i value is called sample quantile and the expected z-score for the point n_i is called the theoretical quantile.

If the data comes from a normal distribution then the theoretical and sample quantiles agree hence giving an approximately straight line Q-Q plot. otherwise, the plot is not a straight line.

* Statistical Tests of Normality :

1) KS - Test (Kolmogorov Smirnov Test) :

Let X_1, X_2, \dots, X_n come from a continuous distribution P . We want to test the following hypothesis:

H_0 : Samples come from P

H_1 : Samples don't come from P

Let F be the CDF of X under H_0 then the empirical distribution (observed) function is:

$$F_{\text{obs}} = \frac{\mathbb{I}(X \leq n)}{\text{Total \# of obs.}}$$

Let F_p be the CDF associated with H_0 . Then, the KS-test statistic is given by

$$D_n = \max \left\{ |F_p(n) - F_{\text{obs}}(n)| \right\}$$

Procedure :

1) Order your dataset. Let the ordered data be $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ so the $F_{\text{obs}} = \frac{\mathbb{I}(X \leq n)}{\text{T. \# of obs}} = \frac{n}{n}$

2) Find F_p for $X_{(i)}$ for each i and calculate the

value at $T_{F_p}(n) - \text{Fars}(n)$ for n values at n and the maximum of these values is D_n .

3) If we use $\alpha = 0.05$ then $D_{c,0.05} = \frac{1.36}{\sqrt{n}}$, H_0 rejected if $D_n > D_{c,0.05}$

2) Shapiro Wilk's Test :

Calculate the W-statistic to check whether x_1, \dots, x_n comes from a normal distribution (KS applies to any dist but SW only applicable for normal).

The W-statistic is given by : $W = \frac{\sum_{i=1}^n (a_i x_{(i)})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$

where $x_{(i)}$ is the ordered sample values and $(a_1, \dots, a_n) = \frac{m^T v^{-1}}{c}$ for $c = \|v^{-1} m\|$ where v is the co-covariance matrix of the ordered sample and $m = (\mathbb{E}(x_{(1)}), \dots, \mathbb{E}(x_{(n)}))^T$ (mean values at the ordered sample).

If p-value < 0.05, we reject $H_0 : x_i$'s come from normal against $H_1 : x_i$'s are not normal.

* Suppose we detected that the data is not normal then we need a Box-Cox transformation.

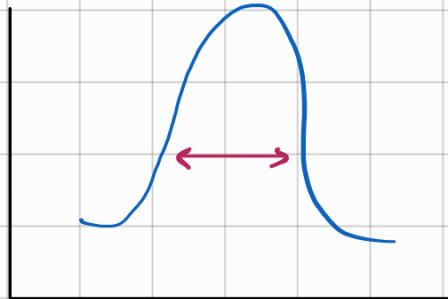
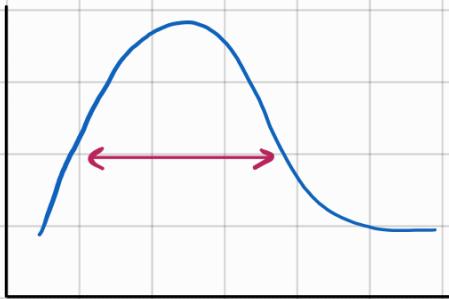
If there is evidence of non-normality then the standard remedy is to transform the response

variables using a Box-Cox transformation.

The response variable in Box-Cox method need to be a +ve value and we assume that there is a transformation parameter λ (hyper parameter) ;

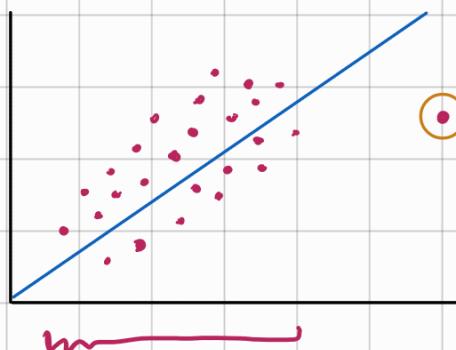
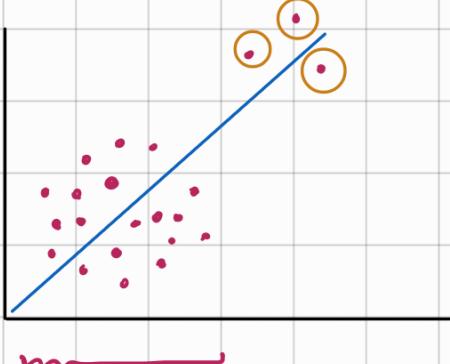
$$Y_i^{(i)} = g(Y_i, \lambda) = \begin{cases} \frac{Y_i^\lambda - 1}{\lambda}, & \lambda \neq 0 \\ \log(Y_i), & \lambda = 0 \end{cases}$$

Assume that the transformed response $f(Y)$ ~ multivariate normal distribution, we try to maximize the likelihood function f^n at response variable wrt λ .



shift x
decrease σ^2

* Outlier : It is a data point whose response Y does not follow the general trend of the data



Leverage obs

Influential obs

If +ve n-coord.
has unusual value
then must at the obs

If it have moderately
unusual n and y coord.

1) Testing for leverage :

$$MLR : Y = X\beta + \varepsilon$$

$$LSE : \hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\text{Fitted Model} : \hat{Y} = X \hat{\beta} = X (X^T X)^{-1} X^T Y = HY$$

H plays an important role in identifying leverage point
point h_{ii} is the i^{th} diagonal element of H.

$$h_{ii} = n_i (X_i^T X_i)^{-1} n_i^T ; \quad X_{n \times K} = \begin{pmatrix} n_1^T \\ n_2^T \\ \vdots \\ n_n^T \end{pmatrix}$$

is the standarized distance at i^{th} obs. from the center at the x-coordinates.

We call a point a leverage point if it has an unusual n-coordinates, high h_{ii} value indicates that the i^{th} obs is a leverage point.

$$\bar{h} = \frac{\sum_{i=1}^n h_{ii}}{n} = \frac{\text{Tr}(H)}{n} = \frac{\text{Rank}(H)}{n} = \frac{K}{n}$$

as H is idempotent

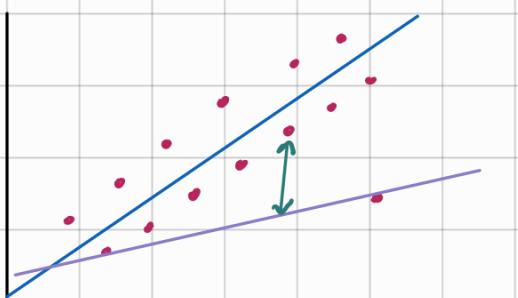
If $h_{ii} > 2\bar{h} = \frac{2K}{n}$ then the i^{th} obs is a possible leverage point

average point.

2) Testing for influential point

Influential : If the point influences any part of a regression analysis, the estimate or slope or regression test results.

The removal of an influential data point causes a large change in the fit especially if the data is small.



There are three ways to determine influential obs.

Cook's Statistic :

For i^{th} obs is based on the difference between the predicted response \hat{Y}_i obtained using the obs. at the predicted response $\hat{Y}_{(i)}$ obtained without the i^{th} obs.

$$D_i = \frac{(\hat{Y}_{(i)} - \hat{Y})^T (\hat{Y}_{(i)} - \hat{Y})}{K \text{ MSRes}} = \frac{\sum_j (\hat{Y}_{(i)j} - \hat{Y}_{ij})^2}{K \text{ MSRes}}$$

squared euclidean distance between the vector of fitted values & the vector of fitted values when i^{th} the obs is deleted

$$D = (D_1, D_2, \dots, D_n)$$

If a value D_i is much larger than others then it may indicate that the i^{th} obs. is influential.

$D_i > 1$: Indicates a potential high influential observation

* DFFITS : Difference b/w the fit statistic

DFFITS indicates the deletion influence at the i^{th} obs. at the fitted value.

For the i^{th} obs., the statistic is given by :

$$DFFITS_i = \frac{\hat{y}_i - \hat{y}_{(i)}}{\sqrt{MSRes(i) h_{ii}}}$$

$\hat{y}_{(i)}$: Fitted value obtained without i^{th} obs

$MSRes(i)$: predicted value at MSRes without i^{th} obs

A possible high value influential obs is indicated by :

$$|DFFITS_i| > 2 \sqrt{\frac{k}{n}}$$

* DFBETAS : How much the regression coeff $\hat{\beta}_j^n$ changes if the i^{th} obs is deleted

$$DFBETAS_{ij} = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\sqrt{MSRes(i)(X^T X)^{-1}}} \quad \forall i=1, \dots, n \quad \forall j=0, \dots, k-1$$

sample coeff.

$\hat{\beta}_{j(i)}$: regression coeff computed without i^{th} obs

A possible high value influential obs is indicated by:

$$|DFBETAS_{ij}| > \frac{2}{\sqrt{n}}$$

$$Y = X\beta + \varepsilon, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\text{Var}(\varepsilon) = \sigma^2 I$$

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$$

When data (y_t, n_t) are collected sequentially in time (time series) the usual assumption of independence of errors is generally not guaranteed such a data is a time series data so $\text{Cov}(\varepsilon_i, \varepsilon_j) \neq 0$ so $\text{Cor} \neq 0$

- * Error autocorrelated / serially correlated correlation b/w errors that are s steps apart

$$\text{Cor}(\varepsilon_t, \varepsilon_{t+s}) = \rho_s, \quad s = 1, 2, \dots$$

- * Correlation b/w residuals is called lagged correlation & it range between $[-1, 1]$

